



GEO: the Gene Expression Omnibus

National Center for Biotechnology Information ■ National Library of Medicine ■ National Institutes of Health ■ Department of Health and Human Services

GEO Database

Examination of gene expression using high-throughput methodologies has become very popular in recent years. Techniques such as microarray hybridization and serial analysis of gene expression (SAGE) allow the simultaneous quantification of tens of thousands of gene transcripts. The Gene Expression Omnibus (GEO) is a public repository that archives and freely distributes high-throughput gene expression data submitted by the scientific community. GEO currently stores approximately half a billion individual gene expression measurements, derived from over 100 organisms, addressing a wide range of biological issues. These huge volumes of data may be effectively explored, queried, and visualized using user-friendly Web-based tools. GEO is accessible at

www.ncbi.nlm.nih.gov/geo

Architecture

Submitters supply their gene expression data in four sections:

Platform: describes the list of elements on the array (cDNAs, oligonucleotides) or the list of elements that may be detected and quantified in that experiment (SAGE tags).

Sample: describes the conditions under which an mRNA source was handled, and the abundance measurement of each element derived from it.

Series: defines a set of related Samples considered to be part of an experiment, and how the samples are related.

Raw data: original microarray scan images or raw quantification data.

Sample data are assembled into biologically meaningful and comparable GEO DataSets (GDS). GDS records provide a coherent synopsis about an experiment and form the basis of GEOs data display and analysis tools.

Submissions

An infrastructure is provided so that submitters can present their data in a MIAME-compliant manner. There are three ways in which data may be deposited with GEO:

Web deposit: Submitters are led through a series of interactive Web forms that accept tab-delimited data tables and accompanying descriptive information.

Direct deposit using Simple Omnibus Format in Text (SOFT): SOFT is a simple line-based format designed for rapid batch submission of data. SOFT files may be readily produced from common spreadsheet applications, and can be uploaded directly to the database.

FTP deposit: Submitters may FTP valid MAGE_ML-formatted reports to GEO.

Submissions are assigned unique and stable GEO accession numbers, and may remain private for several months, typically pending manuscript publication.

Data Mining

The data in GEO can be queried using two NCBI Entrez databases:

Entrez GEO-DataSets provides an **experiment-centric** view of the data in GEO. Experiments of interest may be located by searching for attributes such as free text keywords, technology type, author, organism, and experimental variable information. Once a relevant DataSet is identified, that experiment can be further explored for gene expression profiles of interest (Figure 1) using various supplementary tools provided on the GDS record (Figure 2). Entrez GEO-DataSets is accessible at:

www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=gds

Tools available on the GDS record

- *Cluster heat maps:* A selection of hierarchical and K-means cluster heat maps are provided. Cluster portions of interest can be selected, enlarged, downloaded, plotted as line charts, or linked directly to Entrez GEO-Profiles.
- *Query subset A vs. B:* Using this tool a user can specify that he wants to locate genes displaying 10-fold higher expression values in experimental subgroup A compared to subgroup B, and be directed to Entrez GEO-Profiles matching those criteria.
- *Subset effects:* This feature retrieves all profiles that are flagged as having significant effects with respect to a specific experimental variable, for example 'age' or 'strain'.

Entrez GEO-Profiles provides a **gene-centric** view of the data in GEO. Gene expression profiles (Figure 1) of interest may be located by searching for attributes such as gene name, GenBank accession number, SAGE tag, GDS accession, description, or profiles flagged as having significant effects with regards to specific experimental variables. Entrez GEO-Profiles is accessible at

www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=geo

Tools available within Entrez GEO-Profiles results page

- *Profile neighbors*: returns a list of genes that show a similar expression pattern within a given DataSet.
- *Sequence neighbors*: retrieves profiles related by nucleotide sequence similarity by BLAST
- *Homolog neighbors*: retrieves profiles of genes belonging to the same HomoloGene group
- *Links*: Links to other NCBI Entrez databases including GenBank, PubMed, Gene, UniGene, OMIM, Homologene, Taxonomy, SAGEMap, and MapViewer.

The data in GEO can also be queried outside of the Entrez databases with

GEO BLAST: The GEO BLAST interface allows users to search for GEO-Profiles of interest based on nucleotide sequence similarity. Additionally, all standard BLAST results display 'E' icons that link directly to GEO-Profiles expression data.

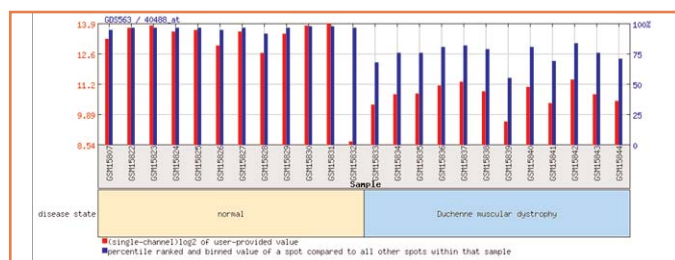


Figure 1: Expression profile of the dystrophin gene in a DataSet examining skeletal muscle biopsies from 12 Duchenne muscular dystrophy patients and 12 unaffected control subjects. Red bars represent gene expression values, blue bars represent intra-sample percentile rank information, providing an indication of the relative expression level of that gene compared to all other genes on the array. Experimental design is reflected in subgroup labels along the bottom of the chart. As expected, the dystrophin gene is seen to be expressed at lower levels in Duchenne patients compared with unaffected control subjects.

Figure 2: GEO DataSet records contain experiment summary information (A), and access to data mining features such as cluster heat maps (B), and a 'Query subset A vs. B' tool that identifies profiles of interest based on experimental variables (C).

Questions relating to GEO submission and GEO query should be sent to: geo@ncbi.nlm.nih.gov

FTP download: All Platform, Sample and Series records, raw data, and GDS value matrices with annotation are available for bulk download via FTP at ftp.ncbi.nih.gov/pub/geo



National Center for Biotechnology Information ■ 8600 Rockville Pike, Building 38A, Room 3S-308
Bethesda, MD 20894 Phone: 301-496-2475 ■ Fax: 301-480-9241 ■ e-mail: info@ncbi.nlm.nih.gov

Web site:
www.ncbi.nlm.nih.gov